

Wright State University

CORE Scholar

Kno.e.sis Publications

The Ohio Center of Excellence in Knowledge-
Enabled Computing (Kno.e.sis)

10-29-2003

What Can Semantics do for Bioinformatics?

Amit P. Sheth

Wright State University - Main Campus, amit@sc.edu

Follow this and additional works at: <https://corescholar.libraries.wright.edu/knoesis>



Part of the [Bioinformatics Commons](#), [Communication Technology and New Media Commons](#), [Databases and Information Systems Commons](#), [OS and Networks Commons](#), and the [Science and Technology Studies Commons](#)

Repository Citation

Sheth, A. P. (2003). What Can Semantics do for Bioinformatics?. .
<https://corescholar.libraries.wright.edu/knoesis/75>

This Presentation is brought to you for free and open access by the The Ohio Center of Excellence in Knowledge-Enabled Computing (Kno.e.sis) at CORE Scholar. It has been accepted for inclusion in Kno.e.sis Publications by an authorized administrator of CORE Scholar. For more information, please contact library-corescholar@wright.edu.



web services

What can Semantics do for Bioinformatics?

Keynote at RCDL 2003, Saint-Petersburg, Russia,
October 29, 2003.

Dr. Amit P Sheth,
Department of Computer Science
LSDIS Lab,
University of Georgia,
Athens,
GA



workflow

proteomics

drug
discovery



Acknowledgements

UGA Biologists/Biochemists

- Will York, Complex Carbohydrate Research Center
- Jonathan Arnold, Fungal Genomics
- Phillip Bowen, CCQC

Project members of LSDIS lab projects
Bioinformatics for Glycan Expression and
METEOR-S (incl. Miller, Kochut, Arpinar)

Special thanks in background research & preparation:
Karthik Gomadam, Christopher Thomas, Kunal Verma

Excellent starting point for complementary material

(a partial list)



- "[Building a Bioinformatics Nation](#)," Lincoln Stein's Keynote at O'Reilly's Bioinformatics Technology Conference 2002
- "[Bio-Ontologies: Their creation and design](#)" Peter Karp, Robert Stevens and Carole Goble
- "[Query Processing with Description Logic Ontologies over Object-Wrapped Databases](#)" Martin Peim, Enrico Franconi, Norman Paton and Carole Goble
- "[Ontologies for molecular biology and bioinformatics](#)" Steffen Schulze-Kremer (paper)
- "Can we do better than Google? Using semantics to explore large heterogeneous knowledge sources," Anatole Gershman, [SWDB Workshop](#), 2003.

Some current BioInformatics Systems



- Tambis
- BioMediator
- Biodas
- BioSem

Data integration, in some cases using an ontology. Single access point for multiple biological information sources; querying multiple sources.



Outline of this talk...

- A Short History of Science
- Challenges in biology
- What can BioInformatics do for Biology?
- What can Semantics do for BioInformatics?
 - Some examples of Semantics-powered Bioinformatics



What is difficult, tedious and
time consuming now ...

What genes do we all have in common?*

Research to answer this question took
scientists two years**

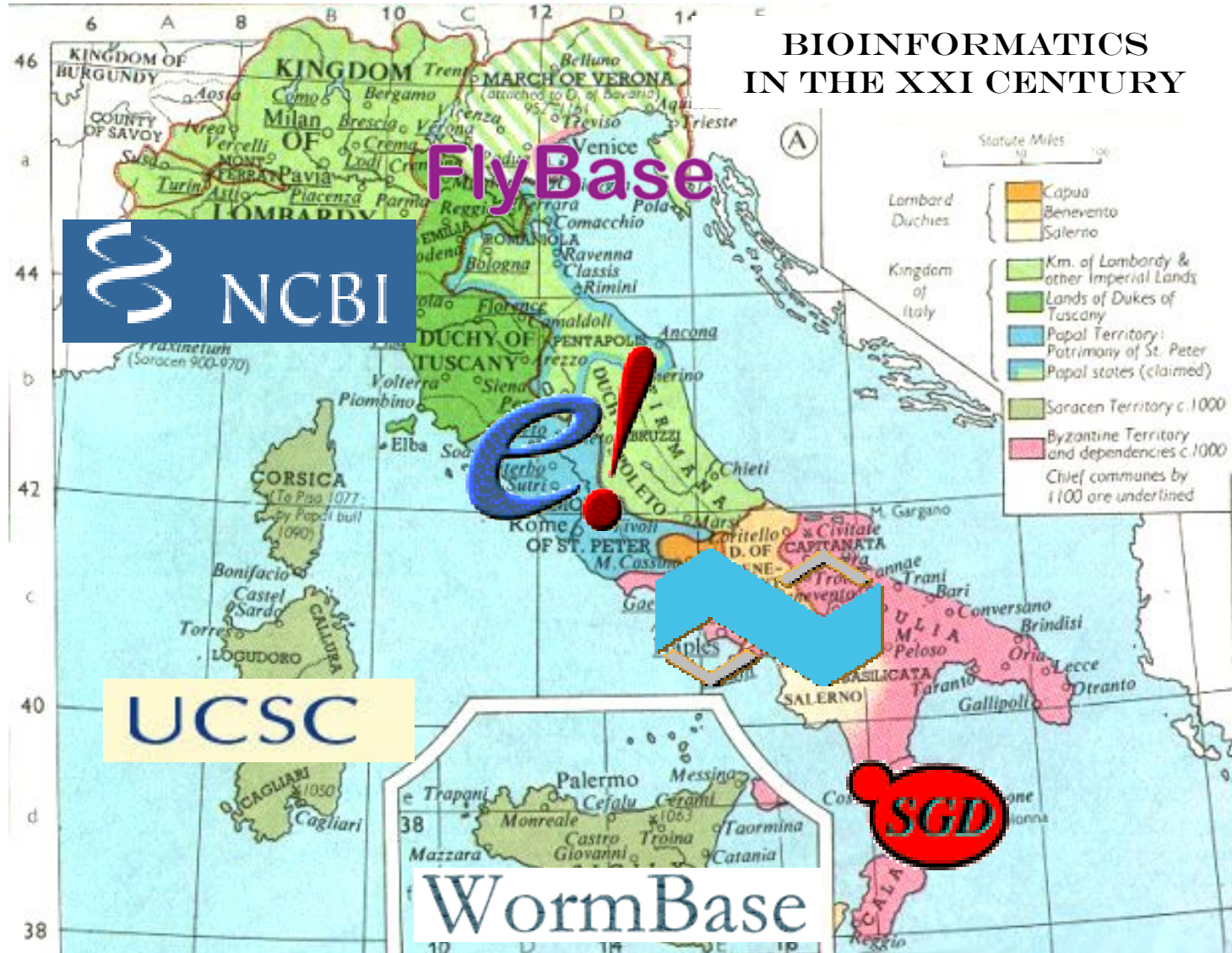
* G. Strobel and J Arnold. Essential Eukaryotic Core,
Evolution (to appear, 2003)

**but we now believe with semantic techniques and
technology, we can answer similar questions much
faster



Why? Bioinformatics, ca. 2002

BIOINFORMATICS IN THE XXI CENTURY



Science then, then and now



In the beginning,
there was thought
and observation.

Science then, then and now

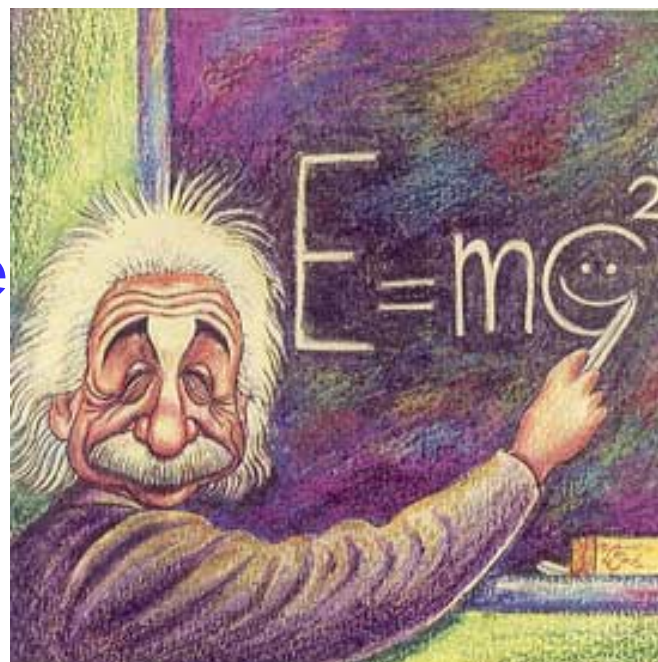
For a long time this didn't change.

- Man thought it would be enough to reason about the existing knowledge to explore everything there is to know.
- Back then, one single person could possess all knowledge in his cultural context.



The achievements are still admirable ...

...as we can see

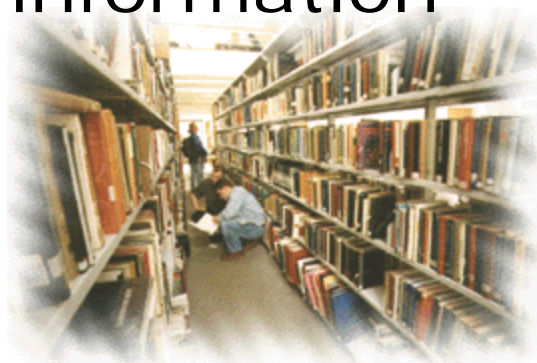


Reasoning and mostly passive observation were the main techniques in scientific research until recently.

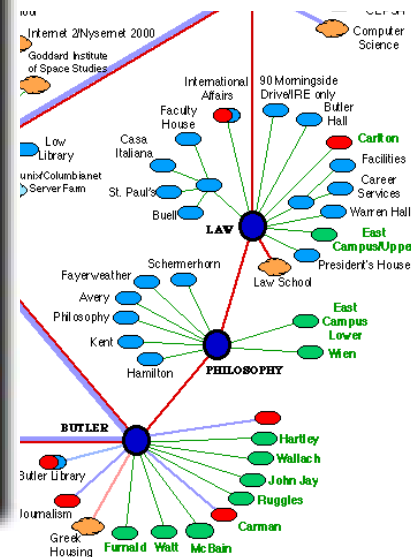
Science then, then and now



A vast
amount of
information



Woodbridge Armstrong



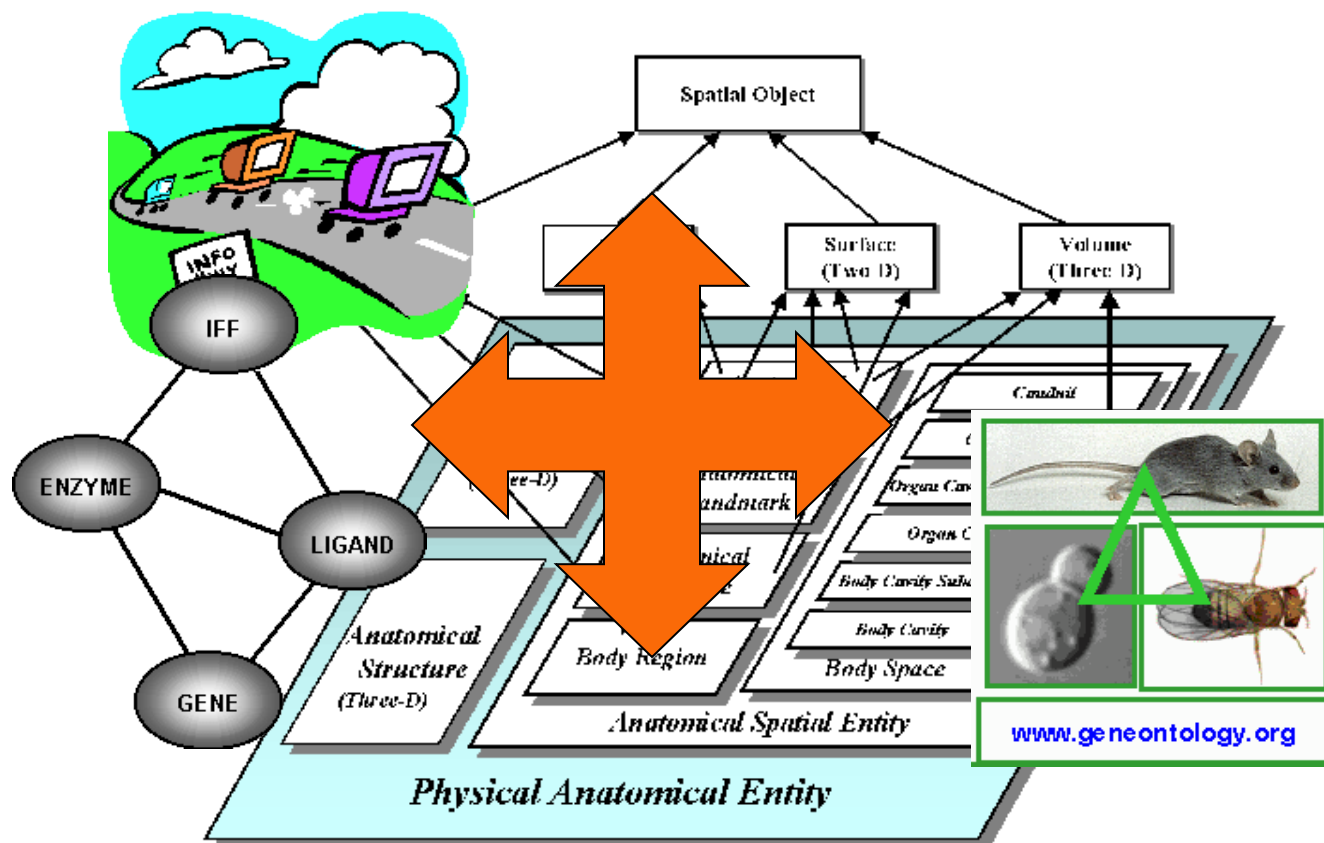
Science then, then and now

No single person,
no group has an
overview of what
is known.

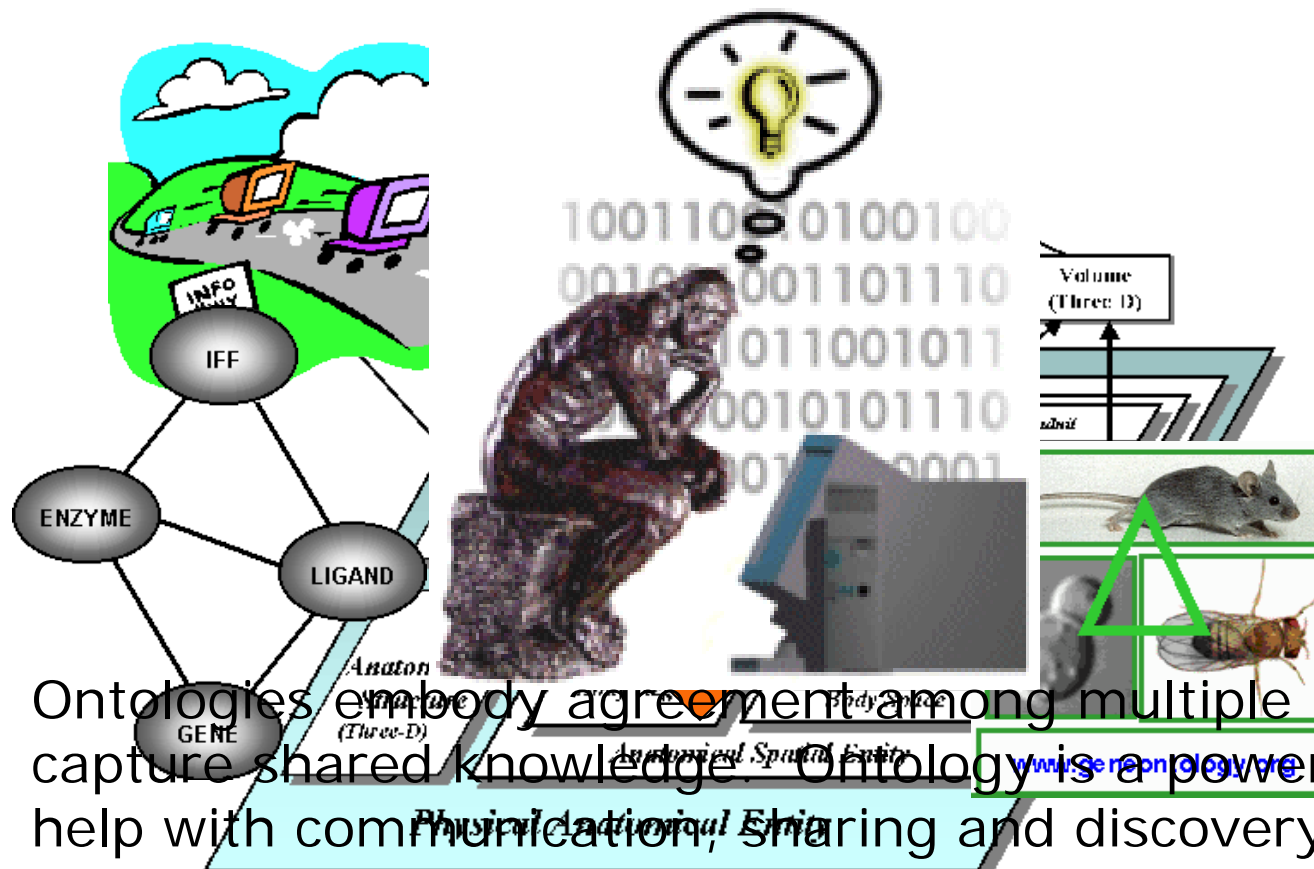


Known, But not known ...→ not known

Science then, then and now



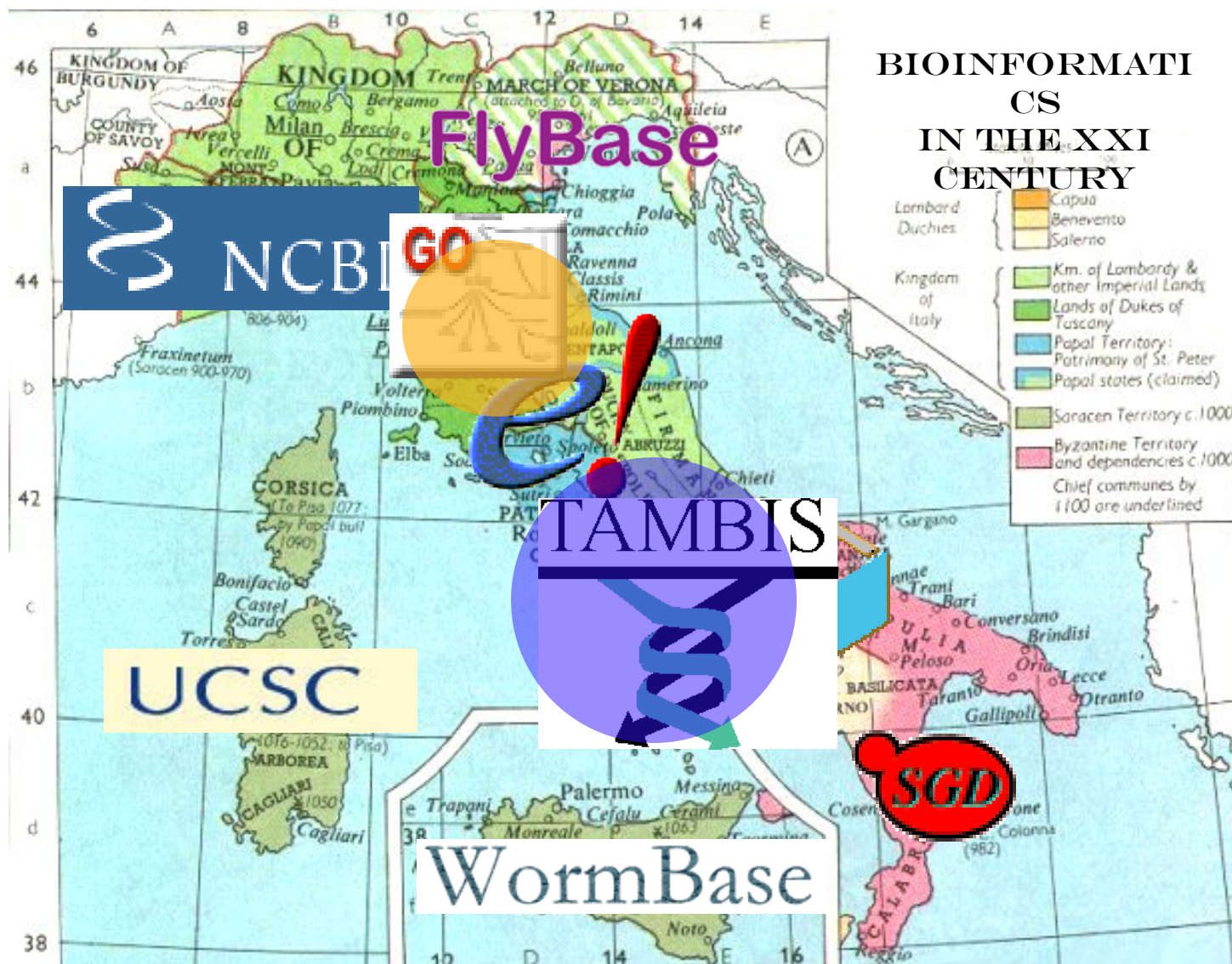
Science then, then and now



Ontologies embody agreement among multiple parties and capture shared knowledge. Ontology is a powerful tool to help with communication, sharing and discovery. We are able to find relevant information ([semantic search/browsing](#)), connect knowledge and information ([semantic normalization/integration](#)), find relationships between pieces of knowledge from different fields ([gain insight, discover knowledge](#))



Intervention by Ontologies ...





Outline of the talk...

- A Short History of Science
- Challenges in biology
- What can BioInformatics do for Biology?
- What can Semantics do for BioInformatics:
 - Some examples of Semantics-powered Bioinformatics



Challenges in biology

- What makes us ill or unwell?
 - Disease identification, disease inducing agents
- What keeps us healthy and makes us live longer?
 - Drug discovery
- Where do we all come from and what are we made of?
 - Genetics and beyond



... and their implications

- Understand biological structures of increasing complexity:
 - Genes (*Genomics*): 1980s
 - Proteins (*Proteomics*): 1990s
 - Complex Carbohydrates (*Glycomics*): 2000s
- Understand biological processes and the roles structures play in them (biosynthesis and biological processes)



Outline

- Evolution of Science
- Challenges in biology
- What can bioInformatics do for Biology?
- What can Semantics do for bioInformatics?
 - Some examples of Semantics-powered Bioinformatics



What can BioInformatics do?

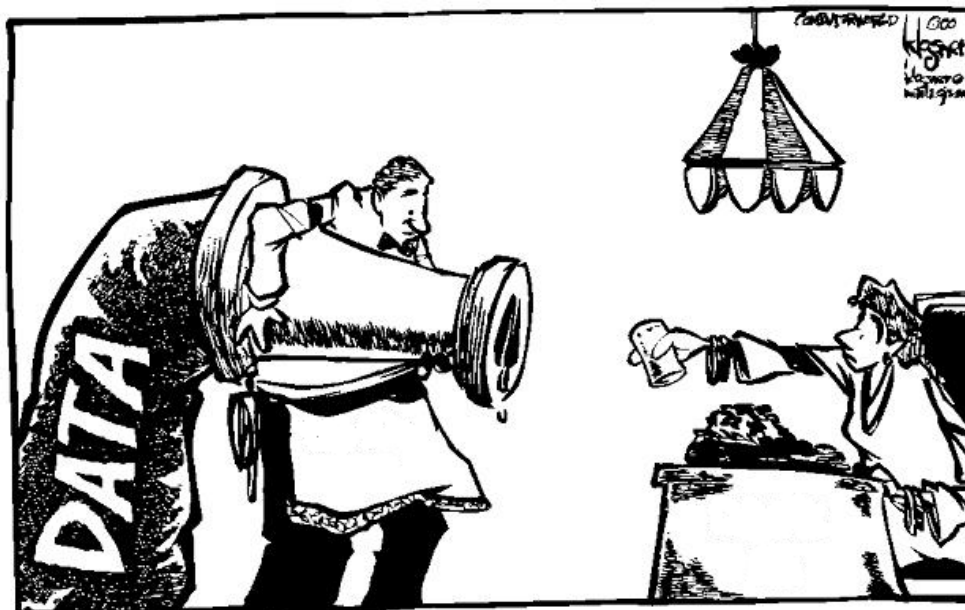
- Analyze genetic and molecular sequences
 - Look for patterns, similarities, matches
 - Identify structures
- Store derived information
 - Large databases of genetic information



Outline

- Evolution of Science
- Challenges in biology
- What can bioInformatics do for Biology?
- What can Semantics do for bioInformatics?
 - Some examples of Semantics-powered Bioinformatics

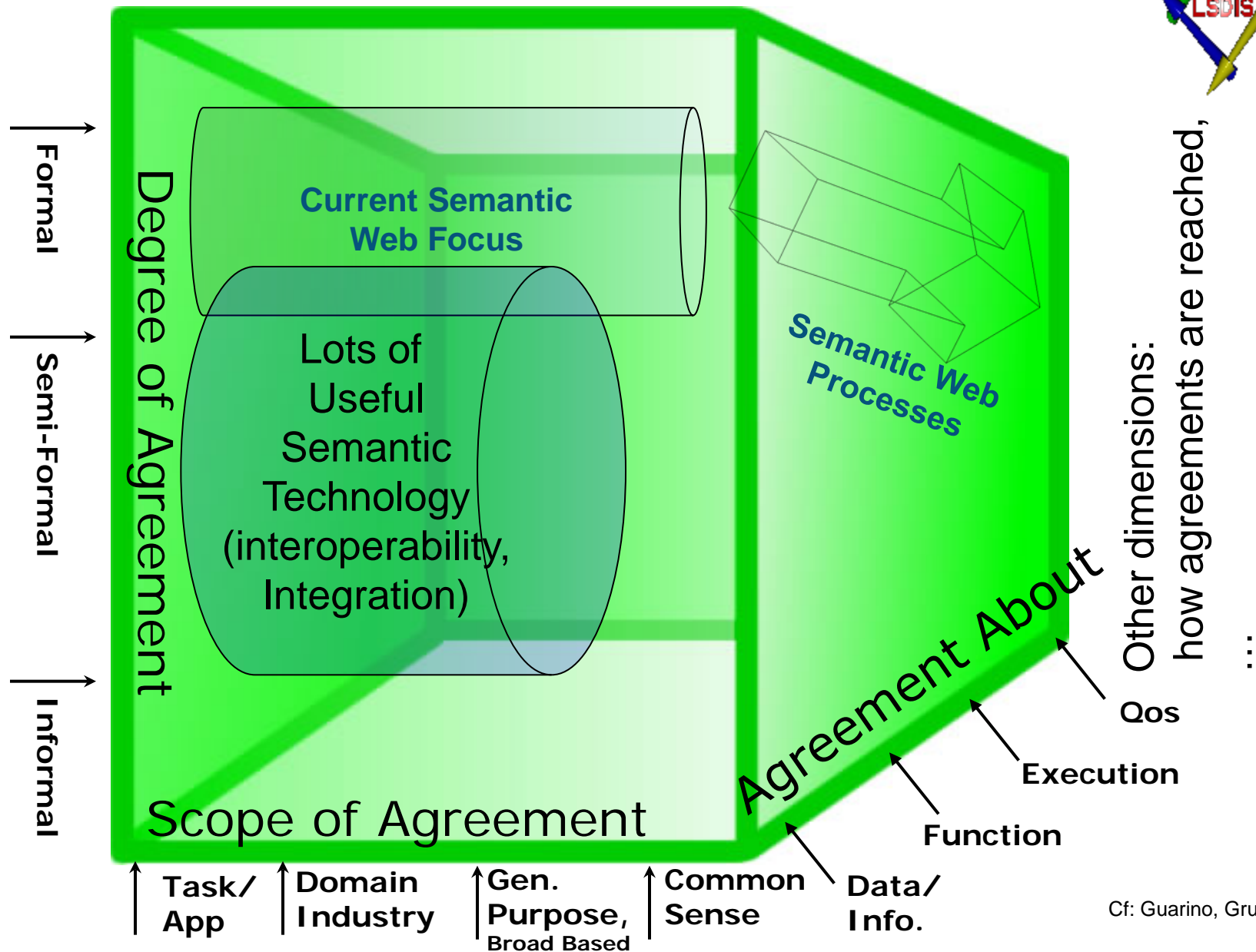
Paradigm shift over time: Syntax -> Semantics



Increasing sophistication in applying semantics & value add

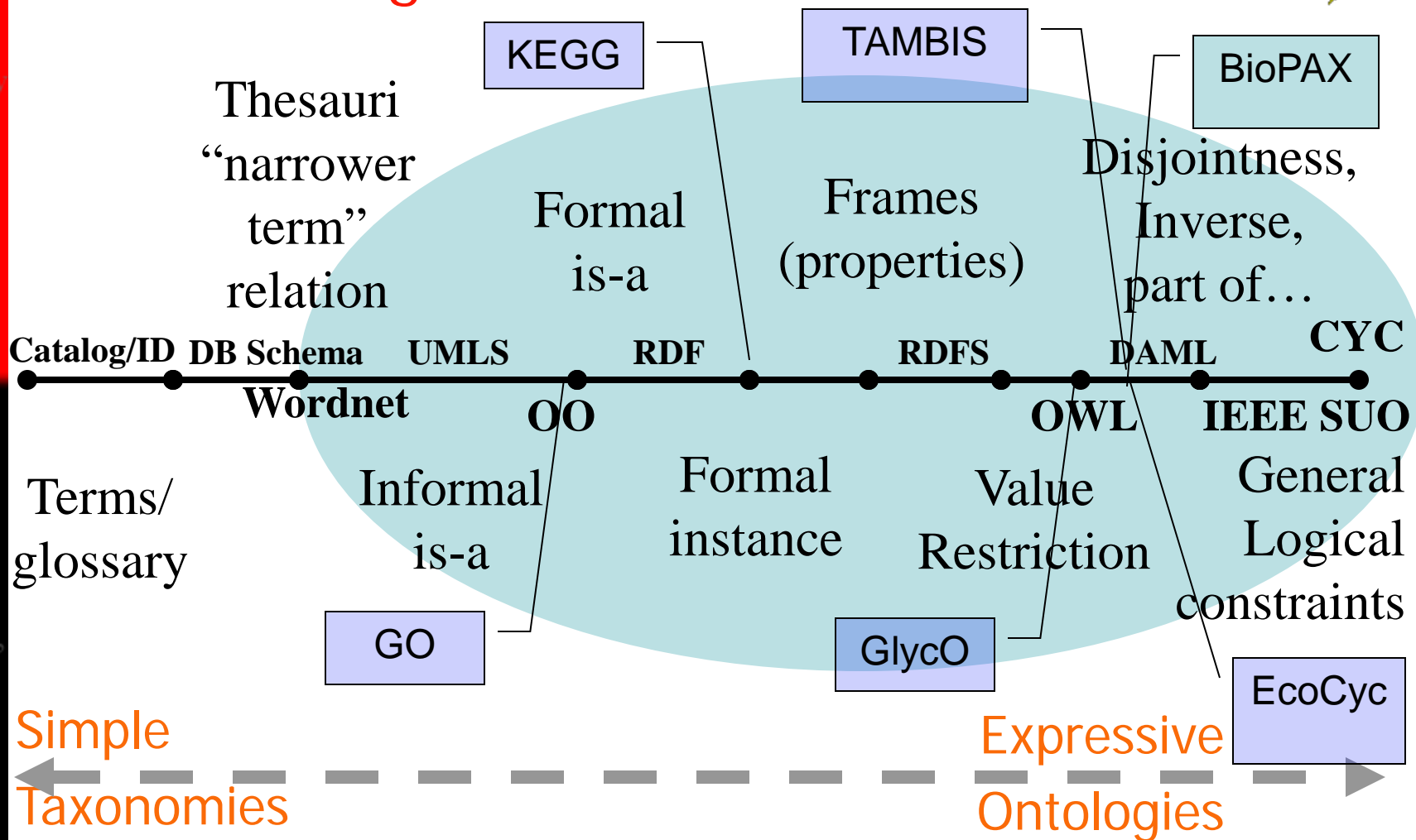
- Relevant Information (Semantic Search & Browsing)
- Semantic Information Interoperability and Integration
- Semantic Correlation/Association, Analysis, Insight and Discovery

Broad Scope of Semantic (Web) Technology



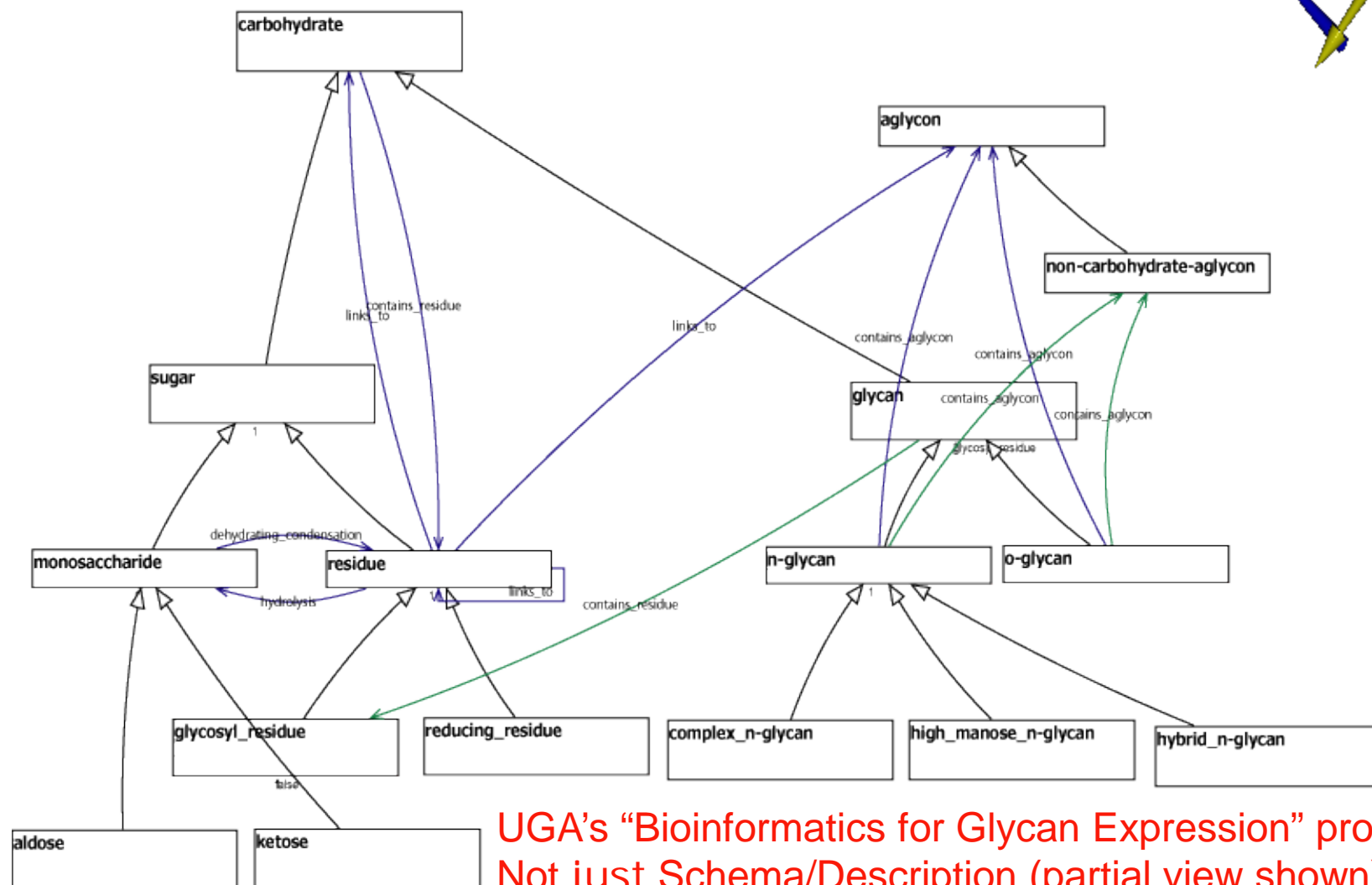
bio
ontologies
b
i
O drug
discovery
i
n L
f S
o D
r I
m S
a t
i
c
s
glycomics

Knowledge Representation and Ontologies





GlycO: Glycan Structure Ontology



UGA's "Bioinformatics for Glycan Expression" proj.
Not just Schema/Description (partial view shown),
also description base/ontology population.
In progress, uses **OWL**.



What can current semantic technology do? (sample)

- Semi-automated (mostly automated) annotation of resources of various, heterogeneous sources (unstructured%, semi-structured, structured data; media content)*
- Creation of large knowledge bases (ontology population) from the trusted sources *
- Unified access to multiple sources*, #
- Inferenceing #
- Relationship/knowledge discovery among the annotated resources and entities; analytics*%
 - Both implicit^ and explicit* relationships

* Commercial: Semagix; %: Near-commercial: IBM/SemTAP;

Commercial: Network Inference; ^ LSDIS-UGA Research



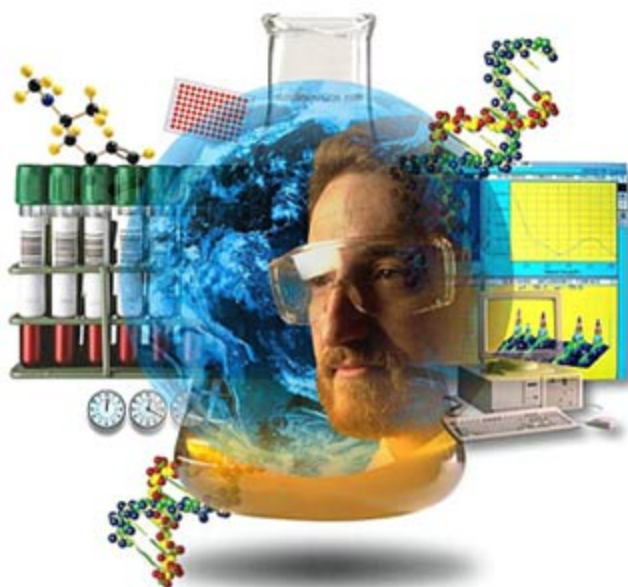
Industry Efforts

(examples with bioinformatics applications only)

- Accenture's Knowledge Discovery Tool (pre-ontology?, not product)
- Semagix's Semantic Browsing and Querying application for drugs for Rogers MIS and its pharmaceutical customers (product and applications); also semantic analysis application (not discussed here)
- Network Inference's Cerebra server: semantic engineering based bioinformatics system aiding drug discovery (product?)

Existing Systems using Semantics for Bioinformatics

FOCUS: SEMANTIC SEARCH AND BROWSING (with nascent work in discovery)



Show: Homo sapiens

Confidence: 6

▼ Articles published in the past 12 months 4 thru 13 of 397

Parkinson's disease over the last 100 years.
Potential nondopaminergic drugs for Parki...
Preclinical versus clinical neuroprotection.
Position of COMT inhibition in the treatmen...
Sleep attacks—facts and fiction: a critical r...
Is the cause of Parkinson's disease enviro...
Strategies to modify levodopa treatment.
Sleep disorders in Parkinson's disease.
Single-photon emission tomograp...
Renaissance of amantadine in the

Recent
Articles

▼ Authors

Corkin S
Cory-Slechta DA
Cote LJ
Cummings JL
Da Prada M
Dagher A
Dahlström A
Damecoul CL
Daniel SE
Davidson MC

Experts

▼ Organizations that published articles 9 thru 18 of 400

Cambridge Centre for Brain Repair, Univer...
Cellular and Clinical Neurobiology Program...
Center for Brain and Cognition, Department ...
Center for Gene Therapy, Tulane University...
Center for Materials of Brain Diseases, Niig...
Center for Molecular and Behavioral Neuro...
Center for Molecular and Behavioral Neur...
Center for Neurodegenerative Disease Res...
Center for Neurologic Diseases, Brigham a...
Center for Neuroscience, North Shore Univ...

Organizations

▼ Parent diseases/phenotypes

Neurodegenerative Diseases
Parkinsonian Disorders

Related
Diseases

▼ Disease/Phenotype

Parkinson Disease

▼ Pathways 1 thru 8 of 8

Alzheimer's disease [Homo sapiens]
Amyotrophic lateral sclerosis (ALS) [Homo ...
Huntington's disease [Homo sapiens]
Oxidative phosphorylation [Homo sapiens]
Parkinson's disease [Homo sapiens]
Prion disease [Homo sapiens]
Proteasome [Homo sapiens]
Ribosome [Homo sapiens]

Metabolic
Pathways

▼ Genes

18 thru 27 of 53

Ly [Homo sapiens]
MAPT [Homo sapiens]
MB [Homo sapiens]
MLLT7 [Homo sapiens]
MTBT1 [Homo sapiens]
MTND1 [Homo sapiens]
ND1 [Homo sapiens]
Ndufv2 [Homo sapiens]
NONO [Homo sapiens]
PARK2 [Homo sapiens]

Genes

▼ Proteins

1 thru 10 of 400

14-3-3 protein tau [Homo sapiens]
2-5 oligoadenylate synthetase 3 [Homo sapi...
Oligoadenylate synthetase 2 [Homo sa...
Oligoadenylate synthetase 3 [Homo sa...
24 kDa subunit of complex I [Homo sapiens]
5-hydroxytryptamine 6 receptor [Homo sapi...
54 kDa nuclear RNA- and DNA-binding prot...
Aconitate hydratase, mitochondrial precur...
AD-015 protein [Homo sapiens]
tor [Homo sapiens]

Proteins

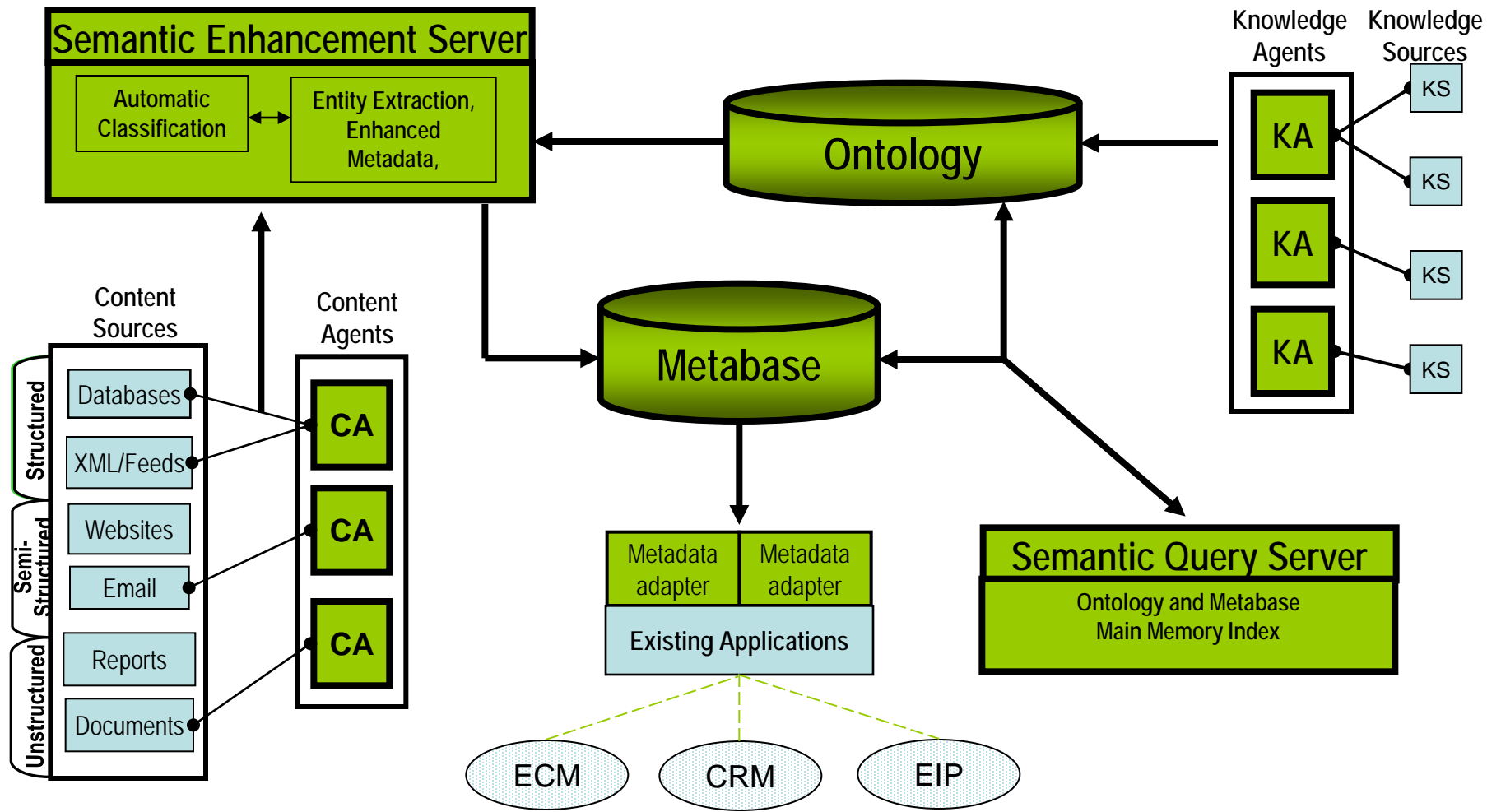
▼ Protein families

Adrenergic
Alpha-synuclein
Amine oxidase
Amine oxidase (flavin-containing)
Aminoglycoside phosphotransferase
Aromatic-L-amino-acid decarboxylase
Arylamine N-acetyltransferase
Arylamine N-acetyltransferase
Aryldialkylphosph...
Arylester

Protein
Families

Semagix Freedom Architecture

(a platform for building ontology-driven information system)

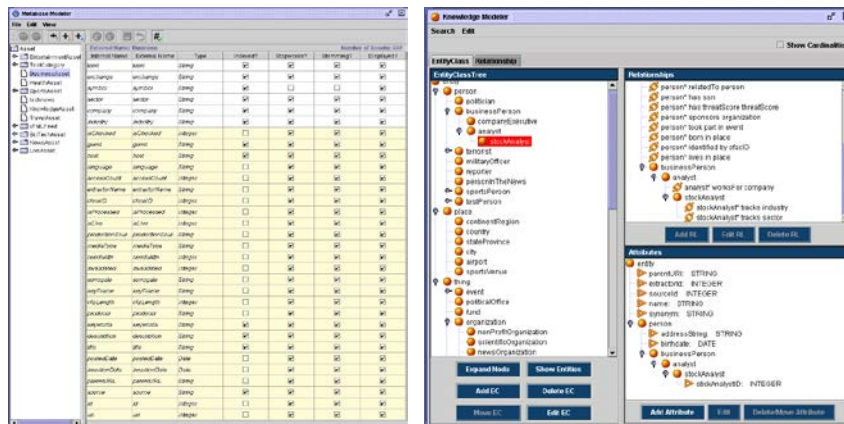


Semagix

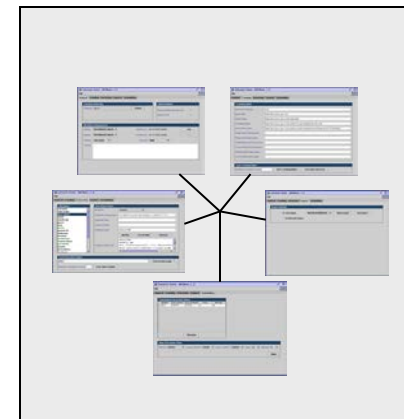
Ontologies Semagix has designed:

- Few classes to many tens of classes and relationships (types); very small number of designers/knowledge experts; descriptive component (schema) designed with GUI
- Hundreds of thousands to over 10 million entities and relationships (instances/assertions/description base)
- Few to tens of knowledge sources; populated mostly automatically by knowledge extractors
- Primary scientific challenges faced: entity ambiguity resolution and data cleanup
- Total effort: few person weeks
- Key requirement: trusted knowledge sources

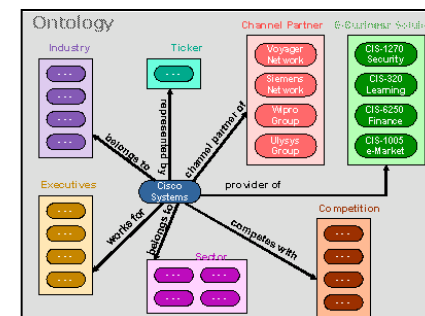
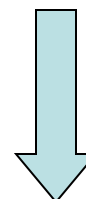
Ontology Creation and Maintenance Steps



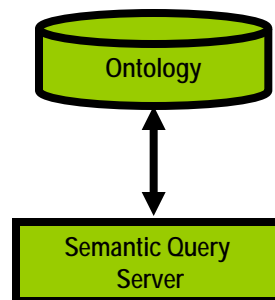
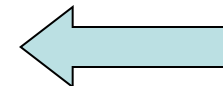
1. Ontology Model Creation (Description)



2. Knowledge Agent Creation



3. Automatic aggregation of Knowledge



4. Querying the Ontology

Help Content Knowledge

Search Results for: drug Procrit in category Medical

FDA notifications. Watch out for counterfeit Procrit, 2 lot
(no description available)

Darbepoetin alfa administered
The objectives of this study were to assess the efficacy and safety of darbepoetin alfa (Aranesp) administered with solid tumors receiving chemotherapy.

Erythropoietic agents as neuroprotectants
Erythropoietin is the primary physiologic regulator of erythropoiesis. It exerts its effect by binding to cell surface receptors. It has been shown that both erythropoietin and its analogs have neuroprotective effects.

Epoetin alfa: current and future indications and nursing implications.
Cancer-related anemia commonly is associated with fatigue and decreased quality of life (QOL). Treatment to achieve hemoglobin targets in patients receiving chemotherapy can improve QOL.

Pure Red-Cell Aplasia and Response to Erythropoietin
To the Editor: Casadevall et al. (Feb. 1, 2000) reported on pure red-cell aplasia and antierythropoietin antibodies in a patient who received recombinant erythropoietin (Eprex).

Role of oral versus IV iron supplementation in the treatment of iron deficiency
BACKGROUND: Preoperative treatment with oral iron (Janssen-Cilag; or PROCIT, Ortho Biotech Products, L.P.) or intravenous iron supplementation increases the erythropoietic response.

Erythropoietin (Procrit; Epogen)
(no description available)

Role of iron in optimizing response to erythropoietin
Approximately 50% of cancer patients have iron deficiency. Iron deficiency was traditionally treated with oral iron. Since the late 1980s, recombinant human erythropoietin (rHuEPO, epoetin alfa [Epogen,

Zoom

procrit Find Entity

Classes Instances

Content Details

Epoetin alfa: current and future indications and nursing implications.

Cancer-related anemia commonly is associated with fatigue and decreased quality of life (QOL). Treatment to achieve hemoglobin targets in patients receiving chemotherapy can improve QOL.

Authors Buchsel, Patricia C Murphy, Barbara J

Side Effects fatigue

Drug Class recombinant hormone

Drugs Epoetin Alfa Procrit

Companies Ortho Biotech Products, L.P.

Hormones erythropoietin

Symptoms fatigue

producer PubMed

Java Applet Window

therapy related)

medica

(chronic disease)

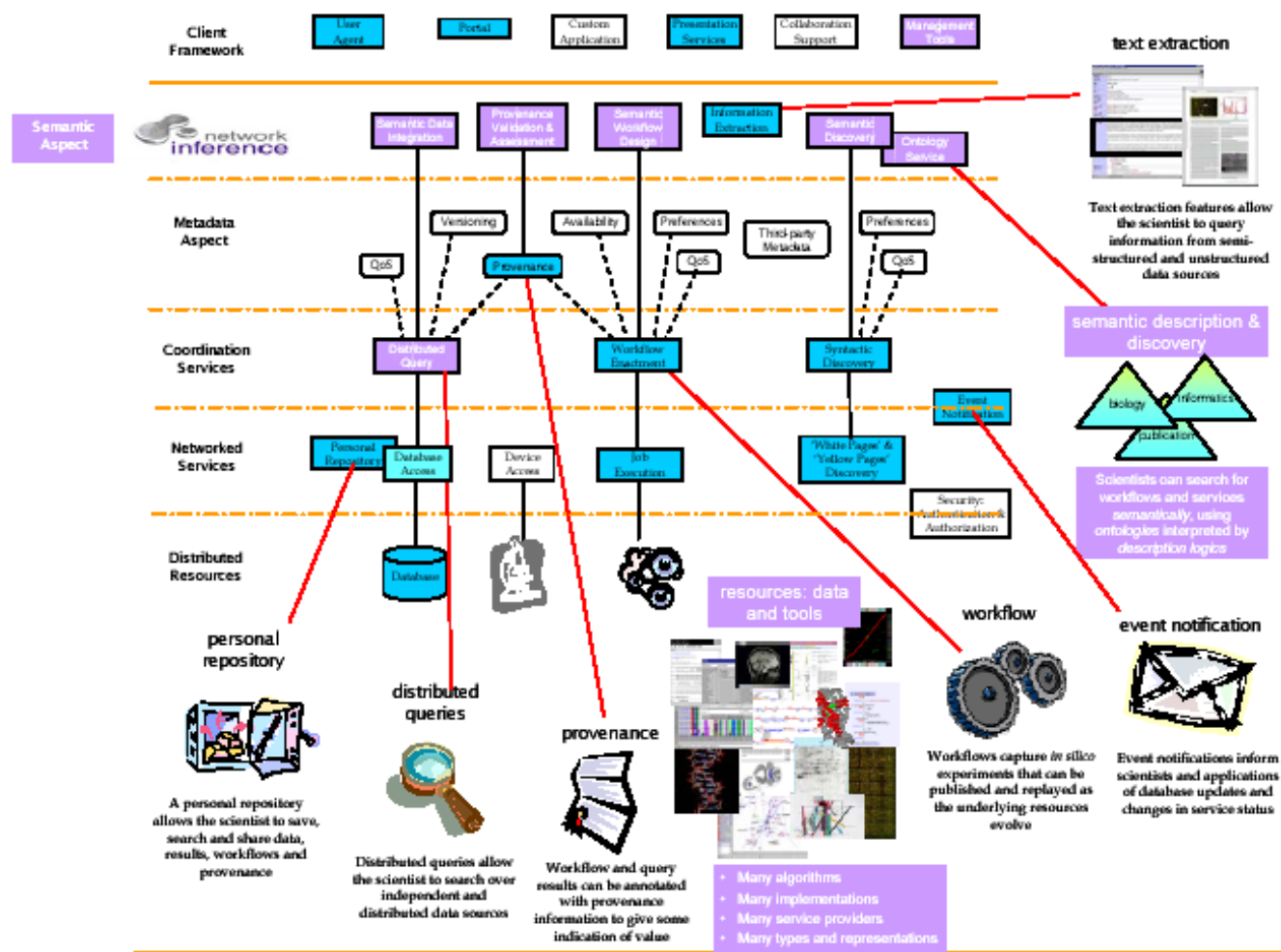
fat

of breath

male appe

blood cell count in

Cerebra's myGrid Framework



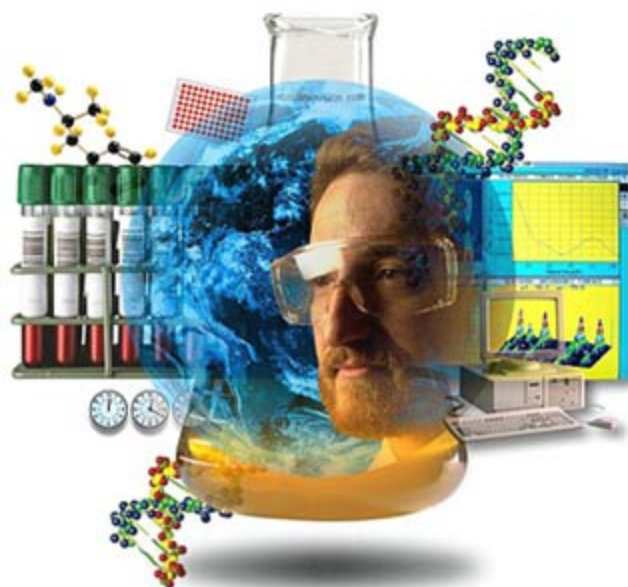


Outline

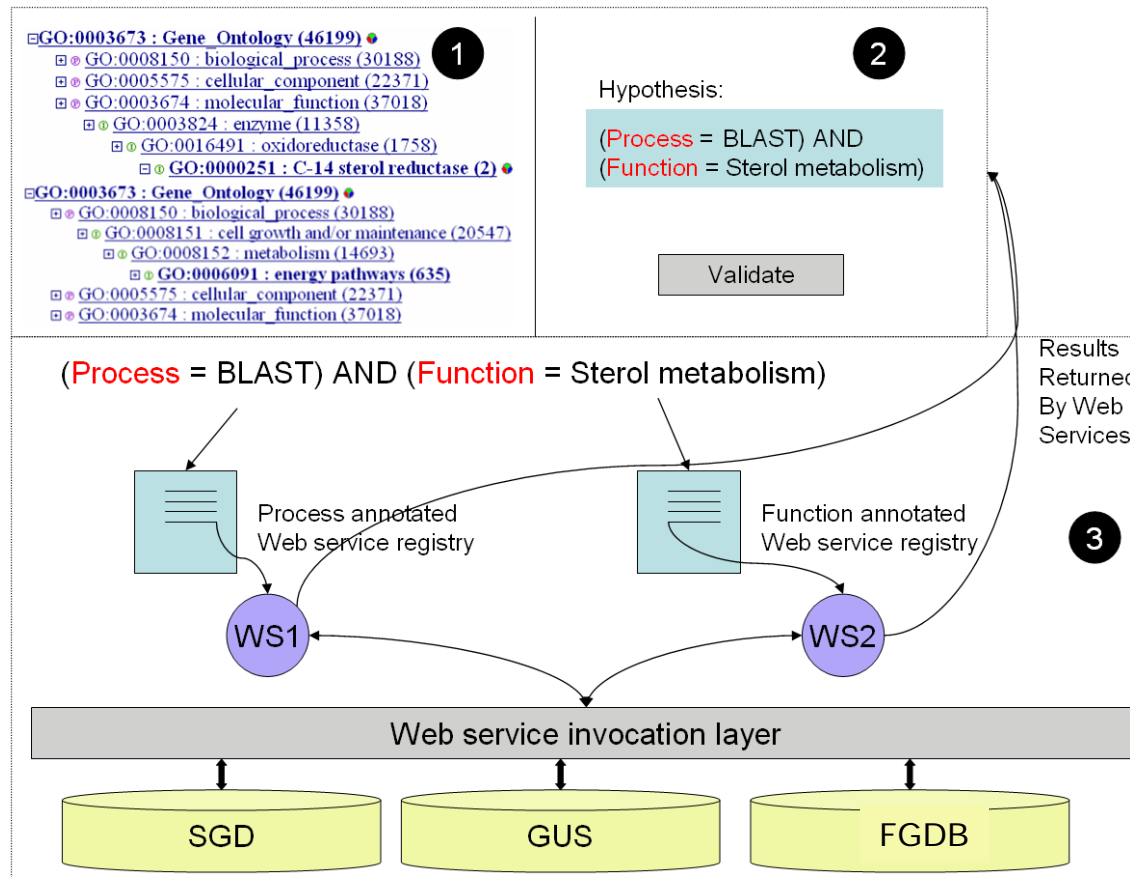
- Evolution of Science
- Challenges in biology
- What can bioInformatics do for Biology?
- What can Semantics do for BioInformatics?
 - Some examples of Semantics-powered Bioinformatics

Applying Semantics to BioInformatics : Example 1

Semantic Browsing, Querying and Integration



Semantic Querying, Browsing, Integration to find potential antifungal drug targets



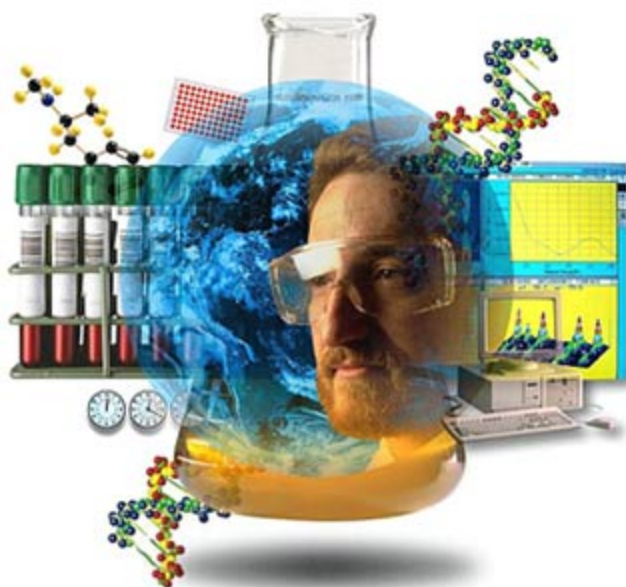
Is this or similar gene in other organism?
(most Antifungals are associated with Sterol mechanism)

Services: BLAST, Co-expression analysis, Phylogeny; If annotated, directly access DB, else use BLAST to normalize

Databases for different organisms

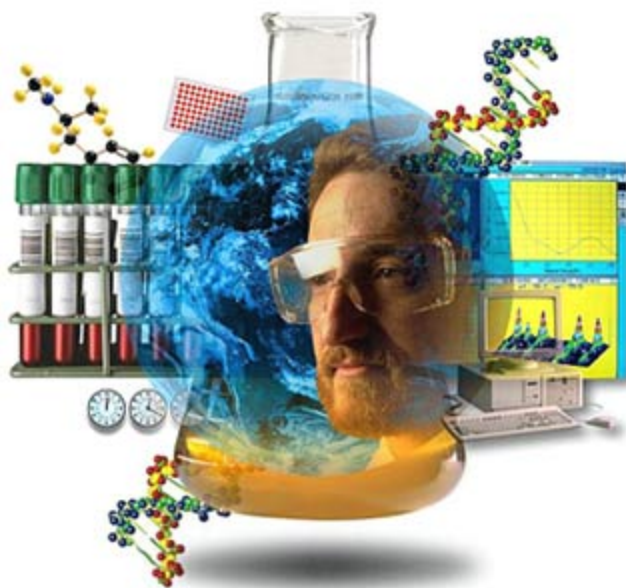
Applying Semantics to BioInformatics : Example 2

Analytics in Drug Discovery



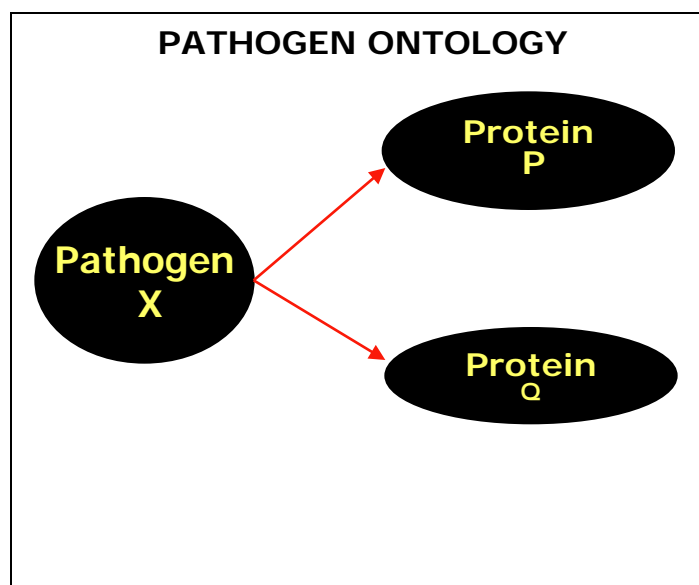
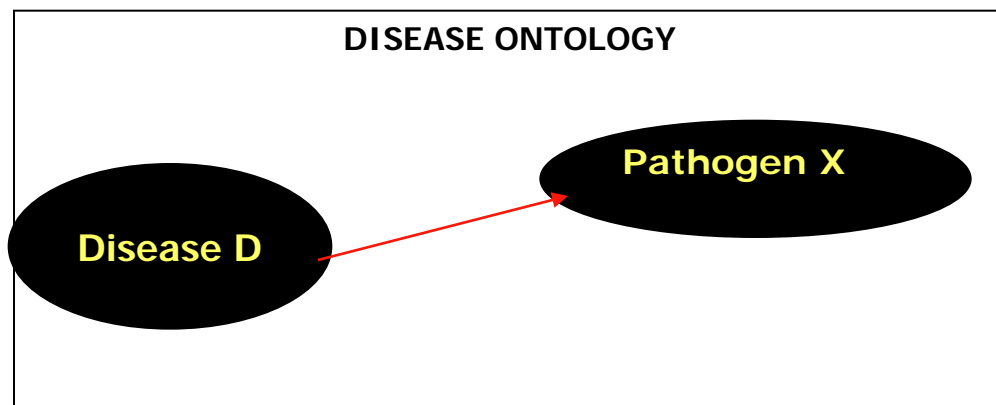
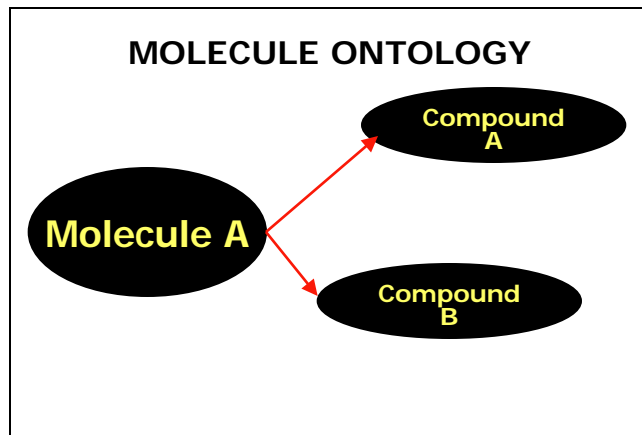
Analytics, Using Explicit and Implicit Relationships in Drug Discovery

- Some molecules contain functional groups that inhibit them from acting as drugs
- Elimination of these groups can make the molecule, a DRUG.





Step 1: Capture domains using ontologies



Step 2: Traverse explicit relationships



STEP 3 (Molecule Ontology):

1. Look up the molecule ontology
2. Identify the composition of the possible drug.

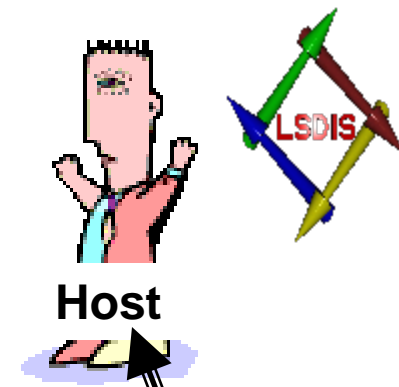
STEP 1 (Disease Ontology):

1. Look up the disease ontology
2. Identify the disease causing pathogen.

STEP 2 (Pathogen Ontology):

1. Look up the pathogen ontology
2. Identify the molecular composition of the pathogen.

Step 3: Discovering Implicit relationships...



Extract the relationships amongst the compounds of the potential drug and the pathogen.

Compound A inhibits the effect of protein P by killing protein P
Compound B inhibits the effect of protein P by killing protein Q

Inferences Based on Relationships



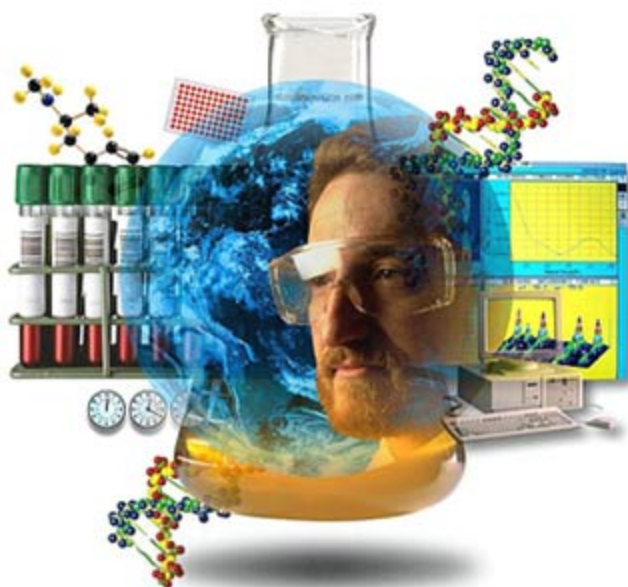
- Compound B doesn't contribute to the curing aspect of the drug, but rather generates a toxin.
- Eliminate compound B and molecule A can be a potential drug.
- However if the host has protein P, then we cannot use protein to bind the drug.
- So look for another drug that can bind at protein Q without producing a toxin
- Eliminate and Discover!!!!



Applying Semantics to BioInformatics : Example 3



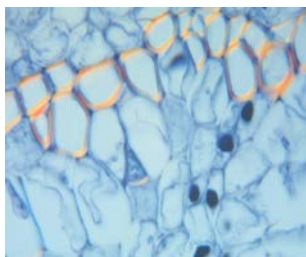
Using Ontologies in cancer research





Disparate Data from Different Experiments

Experiment 1



Metastatic
cancer cells



Increased
GNT-V
Activity

Experiment 2

Cancer marker glycan sequence elevated in glycoprotein beta 1 integrin



Knowledge Stored in Ontologies

- GO Ontology
 - GNT V is an enzyme involved in production of N-glycans
- Glycan Structure Ontology (GlycO)
 - Sequences and structures of Glycans
- Extended Structure Ontology
 - How the structures are made
 - E.g GNT V produces certain class of N-Glycans



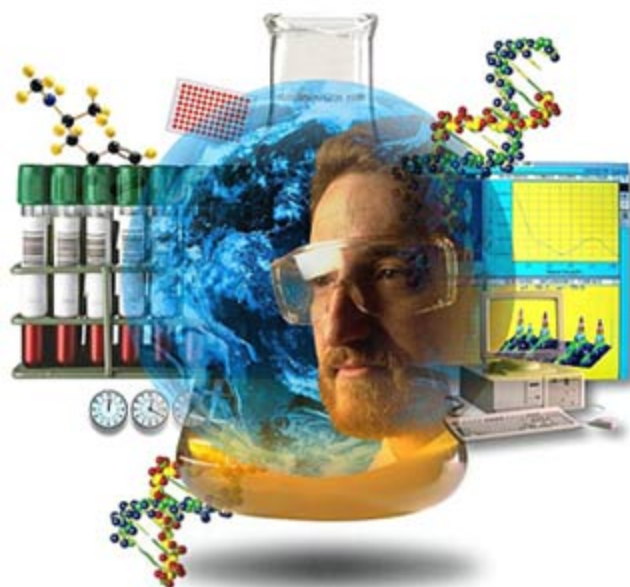
Finding New Information

- Combine data from experiments and knowledge stored in ontologies
- Known assertion from Experiments
 - Beta Integrin is involved in cancer
- New assertion to be tested
 - Are any other glycoproteins involved in cancer ?

Applying Semantics to BioInformatics : Example 4



Applying Semantics to BioInformatics Processes





Creating BioSemantic Processes

- **Question:** *What essential genes do we all have in common?*
- Research process for this using current techniques takes long time (2 years)
G. Strobel and J Arnold. Essential Eukaryotic Core, Evolution (to appear, 2003)
- Let us demonstrate use of Functional, Data, QoS and Execution Semantics to automate the process and reduce time

Creating BioSemantic Processes



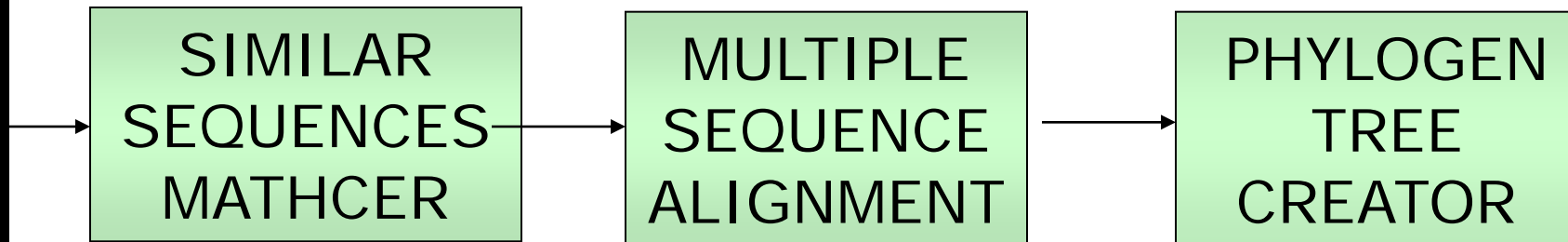
Process for the question:

1. Input GO Id of relevant gene
2. Retrieve similar sequences
3. Perform multiple sequence alignment
4. Construct phylogenetic tree



BioSemantic Process Definition

- Create semantically annotated Web Services (wrapping tools) and process using METEOR-S
- Semantic templates at each node allow choosing multiple tools automatically
- Run multiple instances of process using different tools, combine results

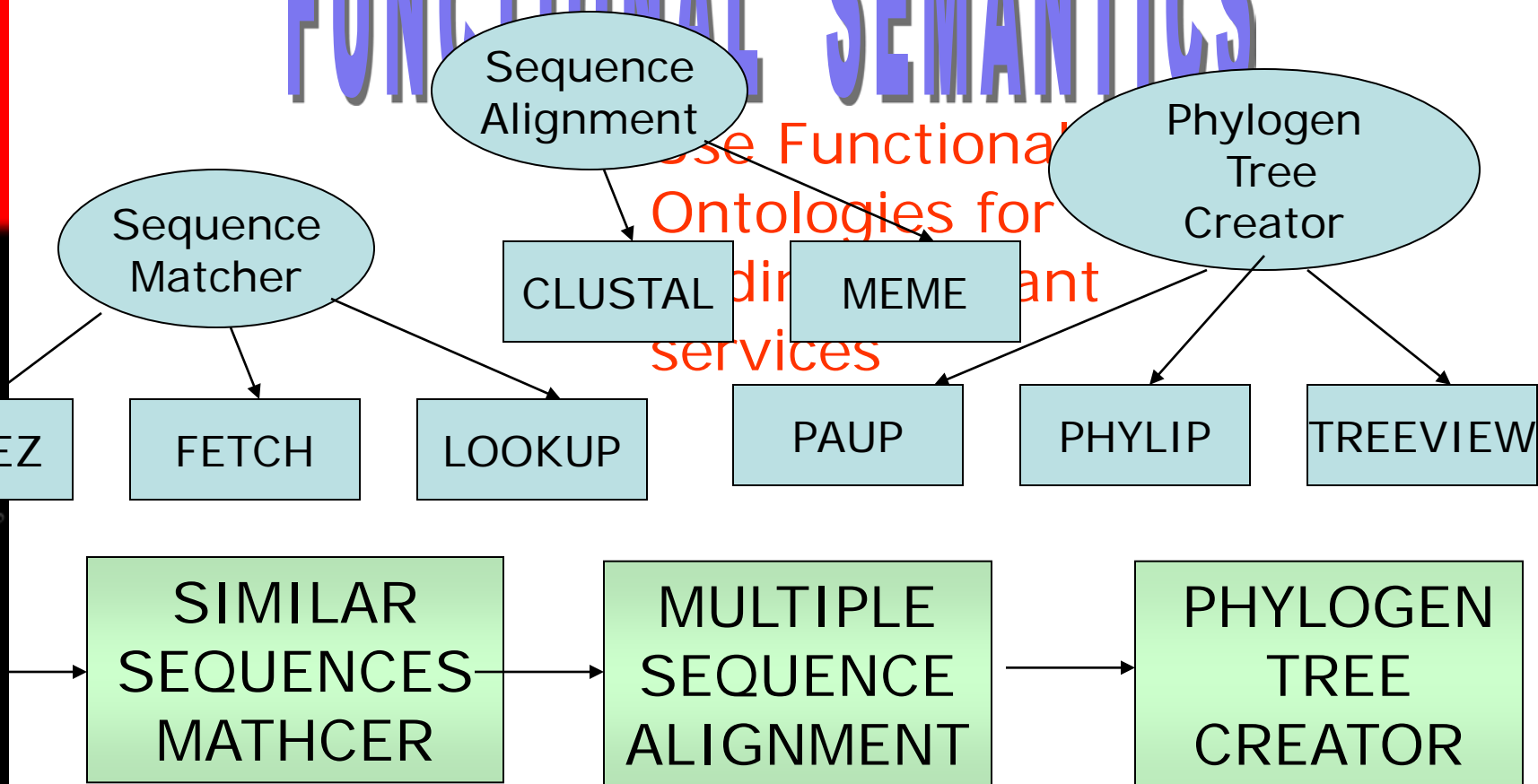




Semantic Bioinformatics Processes

FUNCTIONAL SEMANTICS

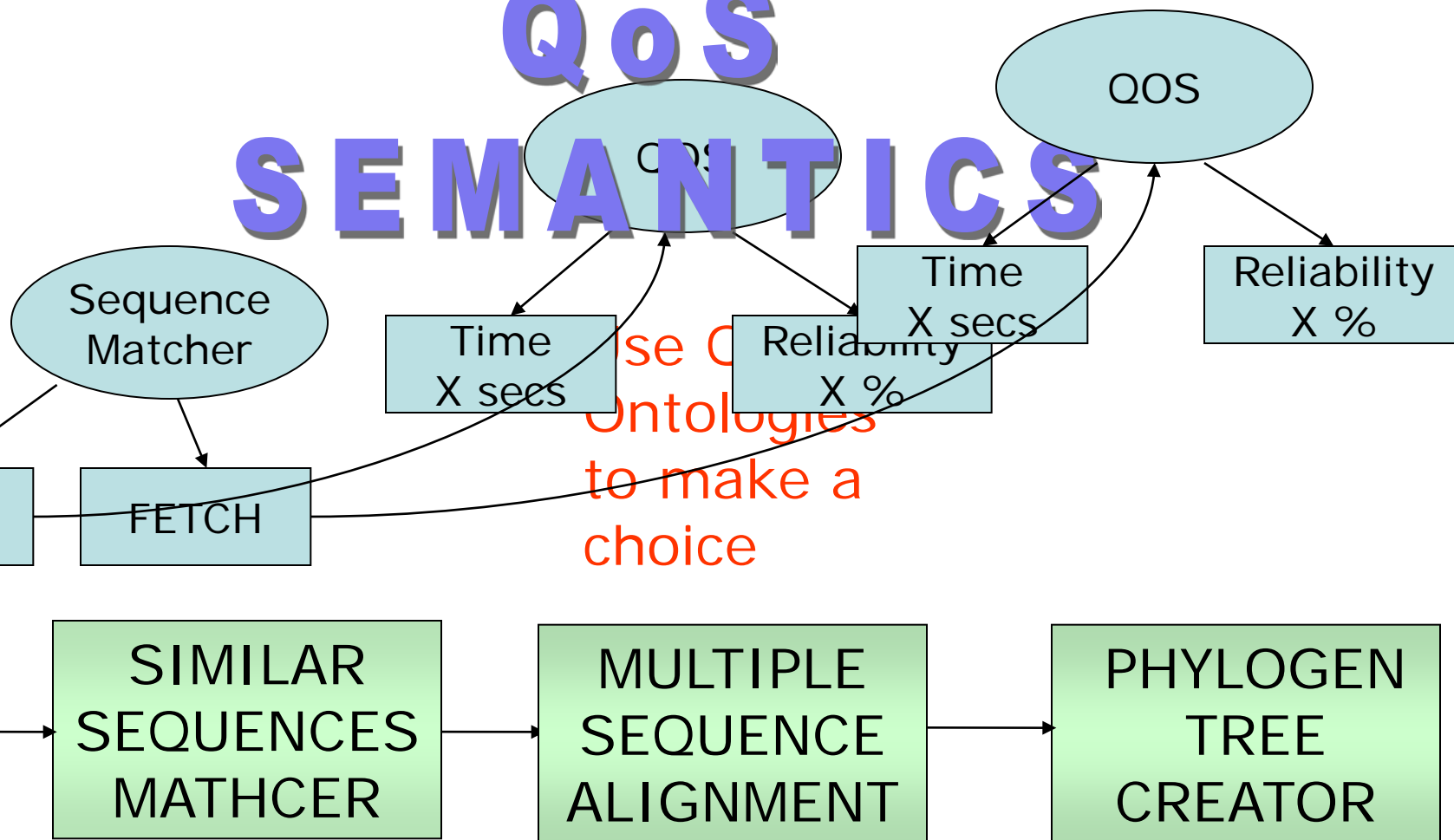
Use Functional
Ontologies for
direct
services





Semantic Bioinformatics Processes

QoS SEMANTICS

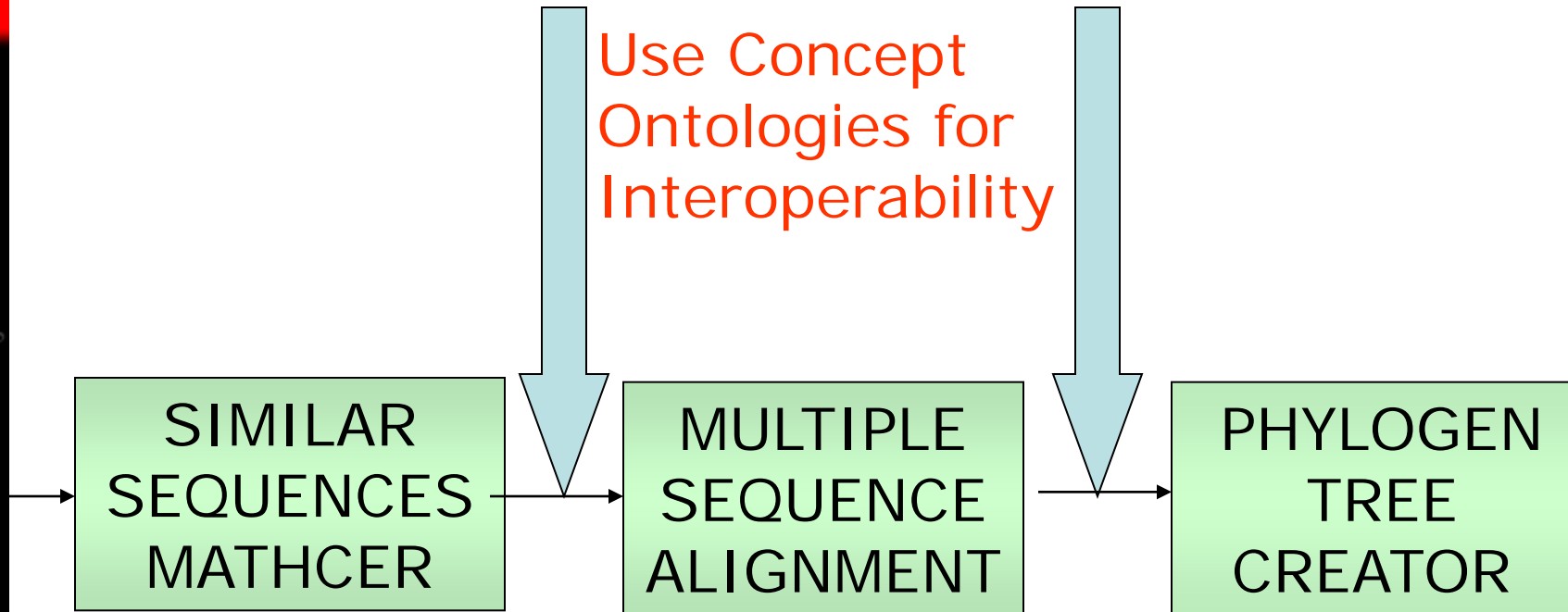




Semantic Bioinformatics Processes

DATA
MAY REQUIRE DATA CONVERSION
SEMANTICS
MAY REQUIRE DATA CONVERSION

Use Concept
Ontologies for
Interoperability

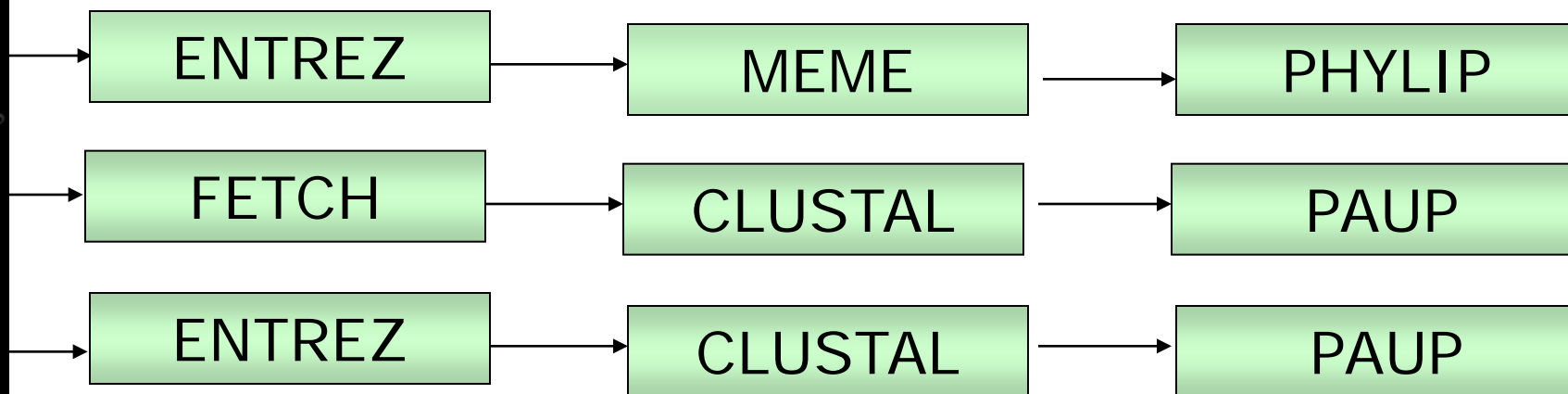




Semantic Bioinformatics Processes

EXECUTION SEMANTICS

Use Execution
Semantics for
execution monitoring
of different instances



Semantic Web Process Mgmt using METEOR-S toolkit



Semantic Web Process Designer

View Process WSDL | View Template | Generate Process | View BPEL Tree | List Ontologies

Control Flow | Data Flow | Process Variables | Service Selection | List Activities

Process Details | Add Web Services | Add Activity Interface | Add Semantic Activity Template | Interface Browser

Activity Name:

Decomposable: ☐

Ontology URL: ▼

Operation Concept:

Discovery URL:

Discovery Specifications:

Ranking Details:

Qos Specifications:

MessagePart Name:

MessagePart Category: ▼

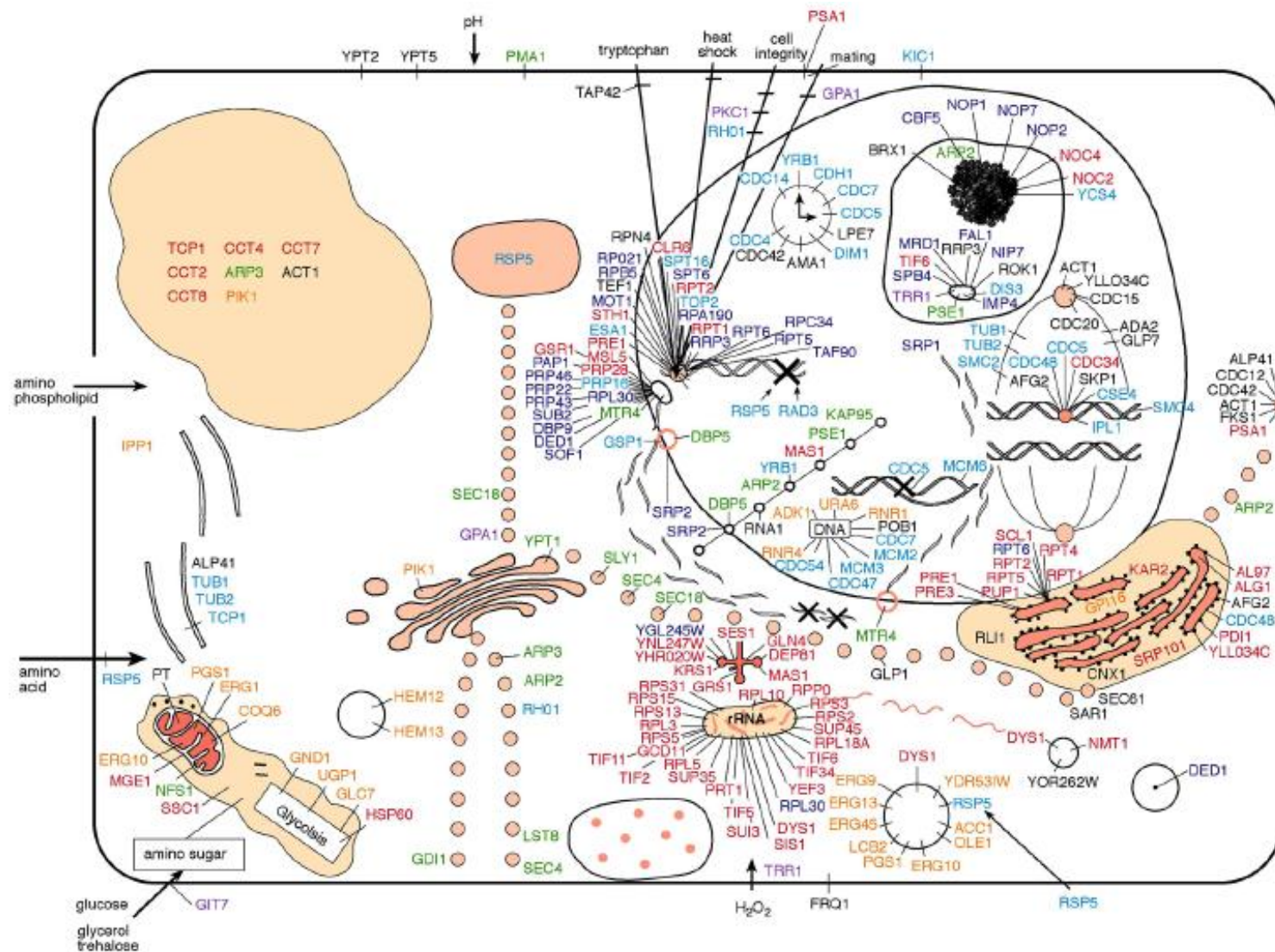
Ontology URL: ▼

Ontological Concept:

MessagePart Type: ▼

Template Construction

Common genes





Conclusion

- Biology research at unique standpoint in history
- Earlier
 - Biology divided into many sub-fields
 - Significant progress in those domains
- Present
 - Searching for the bigger picture
 - Need to combine knowledge from sub-fields
 - Disparate sources, terminologies
- Semantics is the vehicle to get to the answers



Conclusion

- Need to capture the bioinformatics domains using ontologies
 - Not just schema, instances are required
- Presented the use of semantics in
 - Search, Integration
 - Analytics, Knowledge discovery
 - Process Automation
- Collaboration between scientists and computer scientists with semantic techniques and tools make hard things easier



Sources

1. Picture on each sub-section title : <http://www.3rdMill.com>
2. Pictures on the title collage were taken from
 - <http://web.mit.edu/beh.480j/www/covermed.gif> - Glycomics
 - http://bcf.bcm.tmc.edu/proteomics_3.jpg - Proteomics
 - <http://www.microsoft.com/spain/msdn/eventos/presentaciones/images/> - Webservices
 - <http://panda.cs.inf.shizuoka.ac.jp/~cs8060/images/ontology.png> - ontology
 - <http://www.chem.agilent.com/cag/other/workflow-330.gif> - RNA analysis workflow
 - http://www.exploratorium.edu/ti/human_body/dna.html - Genomics
 - Molecular Modelling – Principles and Applications by Andrew R Leach - Drug Discovery

Bioinformatics and Semantic Web at LSDIS Lab, UGA:
<http://lsdis.cs.uga.edu/proj/glycomics/>

Also, Semantic Web Processes: <http://swp.semanticweb.org>